# Time for a Change: Exploring New Avenues for Email Collection Preservation

Smithsonian
*Libraries and Archives*

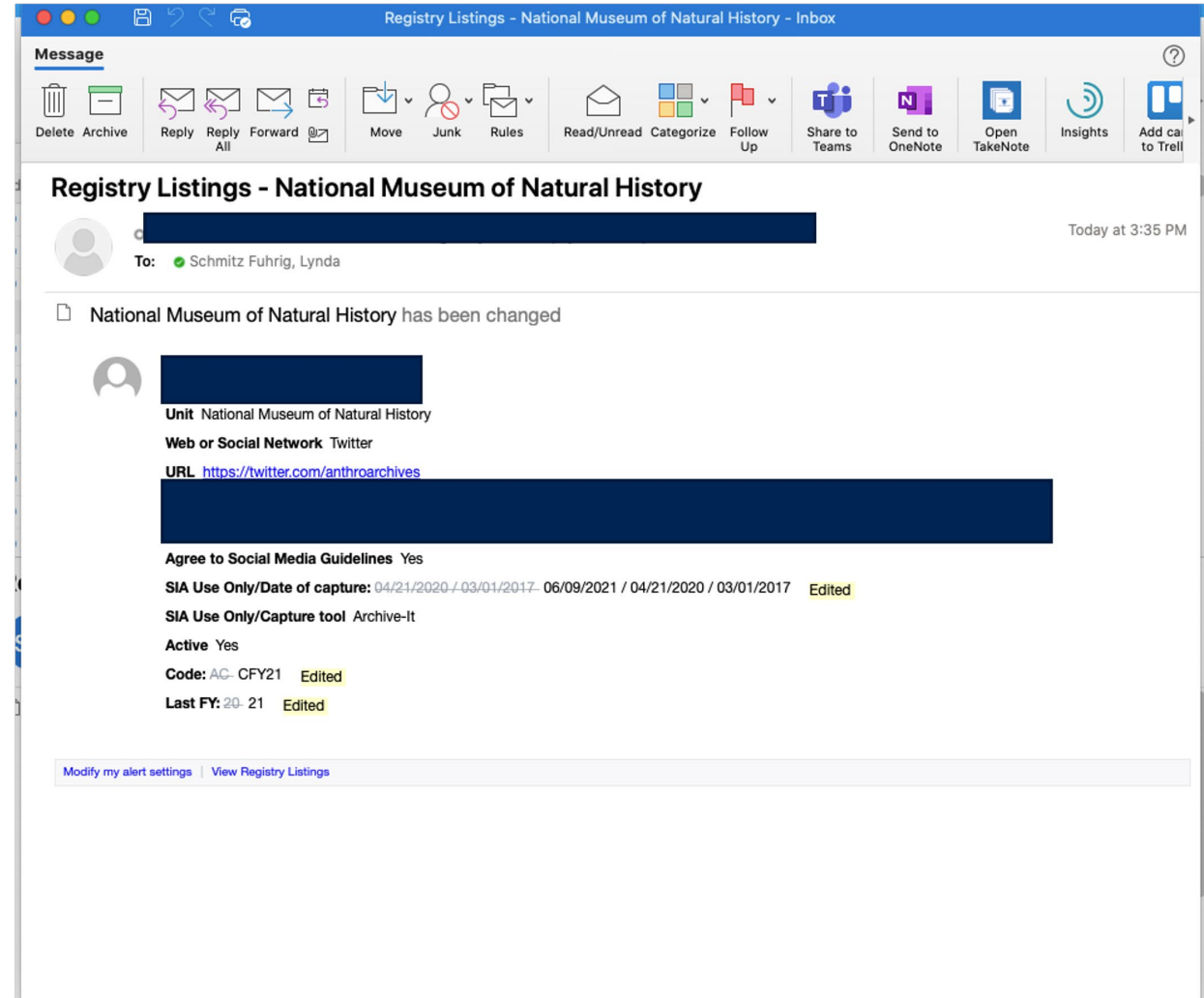Lynda Schmitz Fuhrig

Digital Archivist

Email Archiving Symposium
June 15, 2023

# Through the Years

# Accessioned accounts

- Role at the Smithsonian or specific project – Secretary, undersecretaries, directors, curators, etc.

- When the account holder leaves the Smithsonian, but can be taken in sooner

- No weeding done by SIA except Junk Email or Deleted Items folders in some cases

- 15-year-access restriction. PII permanently restricted

# CERP in 2005

- Smithsonian Institution Archives and Rockefeller Archives Center pilot

- Three Smithsonian offices

- Practices with both email and digital records

- SI migrating from GroupWise to Outlook

- Tested various tools and worked with another project at the State Archives of North Carolina, KY, and PA - EMCAP



Smithsonian
*Libraries and Archives*

# Our Approach



- Open formats

- Anyone could adopt/adapt

- Accounts, not individual emails
  - Easier manageability at this level

# CERP Parser

# Preservation

- Preservation output through Email Account XML Schema (EAXS) using the parser

- Human and machine readable but not "pretty"



Smithsonian
*Libraries and Archives*

# CERP Parser

# DArcMail (Digital Archive for eMail) Suite



Smithsonian
*Libraries and Archives*

# DarcMailXml conversion from MBOX to XML

# Log for conversion to preservation XML



**Convert mbox to XML Summary**

Close

```
########## SETTINGS ##########
account: lsf3
account directory: w:/working/lsf/lsf3
external content: all attachments
chunk size: approximately 1073641824 bytes
output file: w:/working/lsf/lsf3\darcmail_test.xml

########## WARNINGS ##########
skipping duplicate message: <cm.0508451198012.butjykl.qdldjlidd.t@cmail20.com>
in folder darcmail_test (1 times)

########## SUMMARY ##########
total messages in mbox file(s): 6
duplicate messages skipped: 1
total messages in XML output: 5
external content files: 1
```

Smithsonian
*Libraries and Archives*

# DArcMail functions

# Review, Access - DArcMail



Smithsonian
*Libraries and Archives*

# Search results

# Message view

# Export selection to new MBOX file



Smithsonian
*Libraries and Archives*

# Review, Access - ePADD

# 2020s

- Larger accounts
- More people leaving
  = more accessions
- Virus scan tool discontinued

Smithsonian
*Libraries and Archives*

# Email Projects/Events/Reports



Smithsonian Libraries and Archives

# Previous workflow

| Acquire email and create working copy | Scan for viruses | MBOX generation from PST/s | Logs for archival team appraisal | Accession assigned | Preservation processes with XML | Repository Ingest |
|---|---|---|---|---|---|---|

*Bottleneck here due to time to take to complete*

Smithsonian
*Libraries and Archives*

# Current workflow



| Acquire email and create working copy | Scan for viruses | MBOX generation | Logs for archival team appraisal or reviews the PST/s | Accession assigned or declined | MBOX generation | Preservation processes with XML or later | Repository Ingest |

Smithsonian
Libraries and Archives

# Other tools



**libpst Utilities - Version 0.6.76**

## libpst Utilities - Version 0.6.76

## Packages

The various source and binary packages are available at http://www.five-ten-sg.com/libpst/packages/. The most recent documentation is available at http://www.five-ten-sg.com/libpst/devel/.

A Mercurial source code repository for this project is available at http://hg.five-ten-sg.com/libpst/.

This version can now convert both 32 bit Outlook files (pre 2003), and the 64 bit Outlook 2003 pst files. Utilities are supplied to convert email messages to both mbox and for use with many of the CT Summation products. Contacts can be converted to a simple list, to vcard format, or to ldif format for import to an LDAP server.

The libpff project has some excellent documentation of the pst file format.

**Table of Contents**

readpst — convert PST (MS Outlook Personal Folders) files to mbox and other formats
lspst — list PST (MS Outlook Personal Folders) file data
pst2ldif — extract contacts from an MS Outlook .pst file in .ldif format
pst2dii — extract email messages from an MS Outlook .pst file in DII load format
outlook.pst — format of MS Outlook .pst file

# ePADD today – Now accepting PSTs!

# Potential workflow



Acquire email and create working copy

Ingest PST/s into ePADD with Emailchemy tool add in

Review folders for processing in ePADD.

Do weeding of accepted list of general subjects, correspondents, entities, etc. Export new set of MBOX files

Review weeded account in ePADD and/or a log generated from the new ePADD MBOX

If additional processing is done, all Bagged ePADD outputs are ready

Repository Ingest

Smithsonian
Libraries and Archives

- **Don't transfer directories**
  - **Junk**
  - **Personal**
  - **Personnel**
  - **Deleted Items – sometimes ***

*Ingest of PST/s with the Emailchemy for ePADD tool*

All Messages

11.5 GB account

More than 53,000 messages

**6,600**
Correspondents

**2,877**
Entities

**9**
Labels

**45,052**
Image Attachments

**62,285**
All Attachments

**1,237**
Folders

Lexicon Search

Reports

More

# Weeding – do not transfer



Appraisal | Correspondents

All Messages

Search: si email    Show 10 entries

| Name | Sent Messages | Received Messages | Received from Owner |
|---|---|---|---|
| SI Email Announcements | 298 | 4 | 3 |
| SI Email Announcements Prism | 0 | 1 | 0 |

Showing 1 to 2 of 2 entries (filtered from 6,600 total entries)    Previous Next

ePADD Release 10.0.5 © Stanford University

Smithsonian
Libraries and Archives

# Selecting Items



Date: Jun 2, 2010 4:20pm
From: SI Email Announcements <siannounce@si.edu>
To: "SI-GEO-FtPierce, FL" <si-geo-ftpierce@si.edu>,
SI-GEO-National Capital Region <si-geo-ncr@si.edu>, SI-GEO-NYC All Sites <si-geo-nyc@si.edu>,
"SI-GEO-NZP-SCBI Front Royal, VA" <si-geo-nzp-scbi@si.edu>, SI-GEO-Remote Locations <si-geo-remotelocations@si.edu>,
SI-GEO-SAO All Sites <si-geo-sao@si.edu>, "SI-GEO-SERC Edgewater, MD" <si-geo-serc@si.edu>,
SI-GEO-STRI Panama <si-geo-stri@si.edu>
Subject: **Consortia for World Cultures Idea Fair**

Dear Colleagues:

We are moving forward with the development of the ideas or themes that will vitalize the activities and projects of the Consortia as envisioned in the Smithsonian Strategic Plan. One of the suggestions that emerged was to conduct idea fairs that would allow for "bottom-up" emergence of pan-institutional collaborative projects in support of the Consortia.

LABELS ▾    SORT BY ▾    ATTACHMENT VIEW    1 of 302

Do not Transfer

# ePADD Bag



Smithsonian
Libraries and Archives

# First Impressions

- PSTs into ePADD – not complete. Calendar?

- CSVs don't have subject line – but it's coming

- Directory/subdirectory structure not retained on import into ePADD

- Orcids and Sharepoint URLs flagged as CCNs. Finding SSNs



Smithsonian
*Libraries and Archives*

# First Impressions

- Accounts larger than 40 GB have crashed – revisit

- Some formatting issues in the messages. Nature of email though

- Faster way to review and "weed" for final version. ☺

- Appraising archivist – likes the ease of use and the authorities in processing module.



Not enough memory to open this page

Try closing other tabs or programs to free up memory

Error code: Out of memory

Learn more                                   Send feedback

![Smithsonian Libraries and Archives logo]

# Moving ahead

Where do the other previous tools fit in?

New release coming that will have the subject in the CSV. ☺

Tradeoff – willing to live without those Calendar items showing up for appraisal?

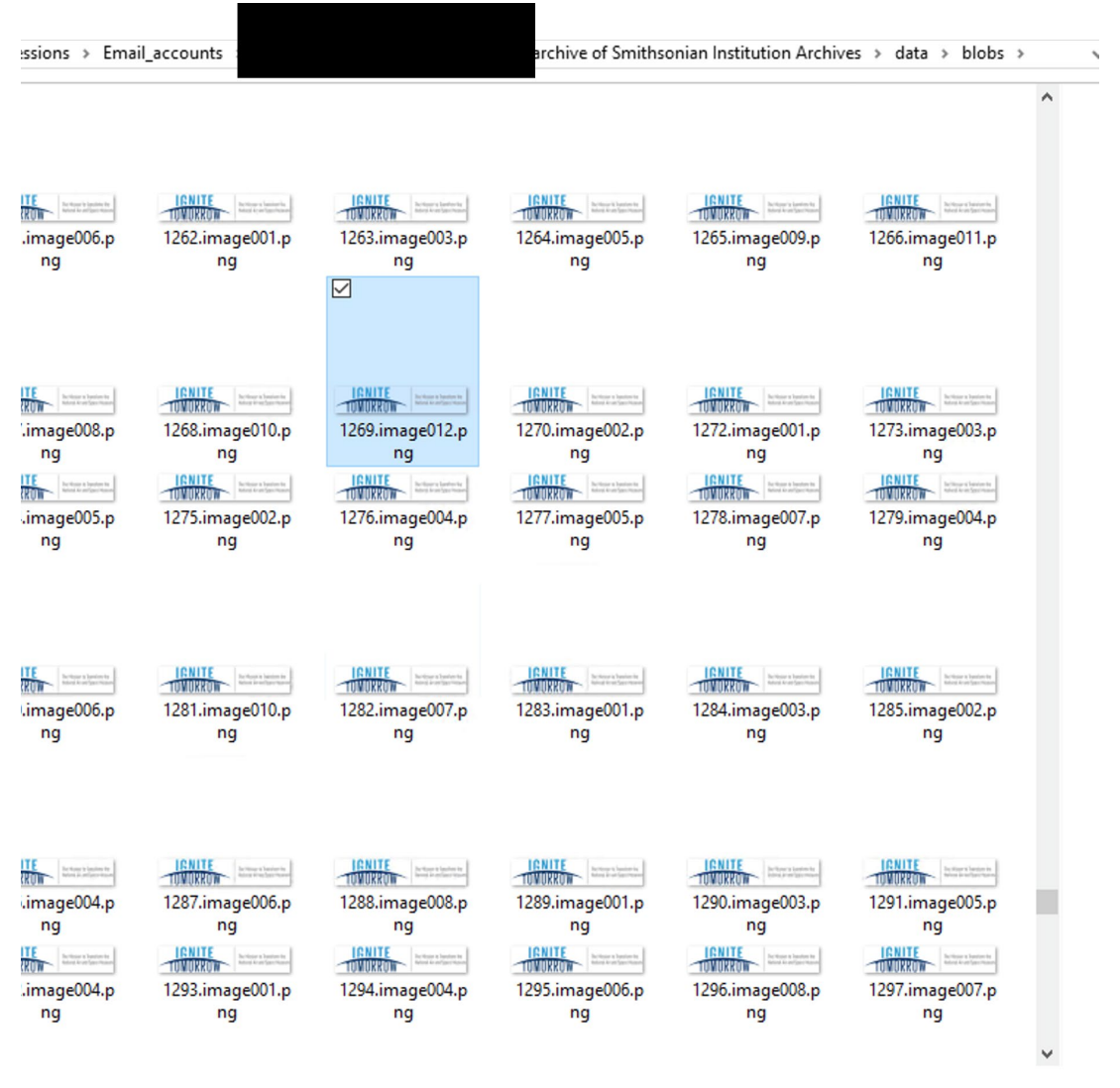How will tasks and calendar items be preserved?

Appraisal – will the archivists want to learn ePADD for reviewing? Or back to the logs?

Smithsonian
*Libraries and Archives*

# Thank you!

## Questions?
[schmitzfuhrigl@si.edu](mailto:schmitzfuhrigl@si.edu)
or @LyndaLSF

Smithsonian
*Libraries and Archives*